

Middlesex University Research Repository

An open access repository of
Middlesex University research

<http://eprints.mdx.ac.uk>

Popescu, Andrei and Traytel, Dmitriy (2019) A formally verified abstract account of Gödel's incompleteness theorems. Fontaine, Pascal, ed. Automated Deduction – CADE 27. CADE 2019. Lecture Notes in Computer Science, vol 11716. In: CADE 27 - 27th International Conference on Automated Deduction, 27-30 Aug 2019, Natel, Brazil. ISBN 9783030294359, e-ISBN 9783030294366. ISSN 0302-9743 [Conference or Workshop Item] (doi:10.1007/978-3-030-29436-6_26)

Final accepted version (with author's formatting)

This version is available at: <https://eprints.mdx.ac.uk/28147/>

Copyright:

Middlesex University Research Repository makes the University's research available electronically.

Copyright and moral rights to this work are retained by the author and/or other copyright owners unless otherwise stated. The work is supplied on the understanding that any use for commercial gain is strictly forbidden. A copy may be downloaded for personal, non-commercial, research or study without prior permission and without charge.

Works, including theses and research projects, may not be reproduced in any format or medium, or extensive quotations taken from them, or their content changed in any way, without first obtaining permission in writing from the copyright holder(s). They may not be sold or exploited commercially in any format or medium without the prior written permission of the copyright holder(s).

Full bibliographic details must be given when referring to, or quoting from full items including the author's name, the title of the work, publication details where relevant (place, publisher, date), pagination, and for theses or dissertations the awarding institution, the degree type awarded, and the date of the award.

If you believe that any material held in the repository infringes copyright law, please contact the Repository Team at Middlesex University via the following email address:

eprints@mdx.ac.uk

The item will be removed from the repository while any claim is being investigated.

See also repository copyright: re-use policy: <http://eprints.mdx.ac.uk/policies.html#copy>

A Formally Verified Abstract Account of Gödel’s Incompleteness Theorems

Andrei Popescu¹ and Dmitriy Traytel²

¹ Department of Computer Science, Middlesex University London, UK

² Institute of Information Security, Department of Computer Science, ETH Zürich, Switzerland

Abstract. We present an abstract development of Gödel’s incompleteness theorems, performed with the help of the Isabelle/HOL theorem prover. We analyze sufficient conditions for the theorems’ applicability to a partially specified logic. In addition to the usual benefits of generality, our abstract perspective enables a comparison between alternative approaches from the literature. These include Rosser’s variation of the first theorem, Jeroslow’s variation of the second theorem, and the Świerczkowski–Paulson semantics-based approach. As part of our framework’s validation, we upgrade Paulson’s Isabelle proof to produce a mechanization of the second theorem that does not assume soundness in the standard model, and in fact does not rely on any notion of model or semantic interpretation.

1 Introduction

Gödel’s incompleteness theorems [10, 13] are landmark results in mathematical logic. Both theorems refer to consistent logical theories that satisfy some assumptions, notably that of “containing enough arithmetic.” The first incompleteness theorem (\mathcal{IT}_1) says that there are sentences that the theory cannot decide (i.e., neither prove nor disprove); the second theorem (\mathcal{IT}_2) says that the theory cannot prove (an internal formulation of) its own consistency. It is generally accepted that \mathcal{IT}_1 and \mathcal{IT}_2 have a wide scope, covering many logics and logical theories. However, when it comes to rigorous presentation, typically these results are only proved for particular, albeit paradigmatic cases, such as theories of arithmetic or hereditarily finite (HF) sets, within classical first-order logic (FOL); and even in these cases the constructions and proofs tend to be “incomplete and (apparently) irretrievably messy” [4, p.16]. Hence, the theorems’ scope remains largely unexplored on a rigorous/formal basis.

The emergence of powerful theorem provers has changed the rules of the game and, we argue, the expectation. Using interactive theorems provers, we can reliably keep track of all the constructions and their properties. Proof automation (often powered by fully automatic provers [18, 28]), makes complete, fully rigorous proofs feasible. And indeed, researchers have successfully met the challenge of mechanizing \mathcal{IT}_1 [15, 25, 27, 35] and recently \mathcal{IT}_2 [27]. Besides reassurance, these verification *tours de force* have brought superior technical insight into the theorems. But they have taken place within the same solitary confinement of scope as the informal proofs.

This paper takes steps towards a more comprehensive prover-backed exploration of the incompleteness theorems, by a detailed analysis of their assumptions. We use Isabelle/HOL [24] to establish general conditions under which the theorems apply to

a partially specified logic. Our formalization is publicly available [31]. An extended technical report gives more details [30].

We start with a notion of logic (Section 2) whose terms, formulas and provability relation are kept abstract (Section 2.1). In particular, substitution and free variables are not defined, but axiomatized by some general properties. On top of this logic substratum, we consider an arithmetic substratum, consisting of a set of closed terms called *numerals* and an order-like relation (Section 2.2). Also factored in our abstract framework are encodings of formulas and proofs into numerals, the representability of various functions and relations as terms or formulas (Section 2.3), variations of the Hilbert-Bernays-Löb derivability conditions [16, 23] (Section 2.4), and standard models (Section 2.5).

Overall, our assumptions capture the notion of “containing enough arithmetics” in a general and flexible way. It is general because only few assumptions are made about the exact nature of formulas and numerals. It is flexible because different versions of the incompleteness theorems consider their own “amount of arithmetics” that makes it “enough,” as proper subsets of these assumptions. Indeed, our formalization of the theorems (Section 3) proceeds in an austere-buffet style: Every result picks just enough infrastructure needed for it to hold—ranging from diagonalization which requires very little (Section 3.1) to Rosser’s version of \mathcal{IT}_1 which is quite demanding. This approach caters for a sharp comparison between different formulations of the theorems, highlighting their trade-offs: Gödel’s original formulation of \mathcal{IT}_1 versus Rosser’s improvement (Section 3.2), proof-theoretic versus semantic versions of \mathcal{IT}_1 (Section 3.2), and Gödel’s original formulation of the \mathcal{IT}_2 versus Jeroslow’s improvement (Section 3.3).

Abstractness is our development’s main strength, but also a potential weakness: Are our hypotheses reasonable? Are they consistent? These questions particularly concern our axiomatization of free variables and substitution—a notoriously error-prone area. As a remedy, we instantiate our framework to Paulson’s semantics-based \mathcal{IT}_1 and \mathcal{IT}_2 for HF set theory [27], also performing an upgrade of Paulson’s \mathcal{IT}_2 to a more general and standard formulation: for consistent (not necessarily sound) theories (Section 4). In the rest of this section, we discuss some formalization principles and related work.

Formal Design Principles Our long-term goal is a framework that makes it easy to instantiate the incompleteness theorems and related results to different logics. This is a daunting task, especially for \mathcal{IT}_2 , where a lot of seemingly logic-specific technicalities are required to even formulate the theorem. The challenge is to push as much as possible of the technical constructions and lemmas to a largely logic-independent layer.

To this end, we strive to make minimal assumptions in terms of structure and properties when inferring the results—we will call this the *Economy* principle. For example, we do not define, but axiomatize syntax in terms of a minimalistic infrastructure. We assume a generic single-point substitution, then define simultaneous substitution and infer its properties. This is laborious, but worthwhile: Any logic that provides a single-point substitution satisfying our assumptions gets the simultaneous substitution for free.

As another instance of Economy, when faced with two different ways of formulating a theorem’s conclusion we prefer the one that is *stronger under fewer assumptions*. (And dually, we prefer weakness for a theorem’s assumptions.) For example, we discuss two variants of consistency: (1) “does not prove false” or (2) “there exists no formula such that itself and its negation are provable” (Section 3.3). While the statements are

equivalent at the meta-level, their representations as object-logic formulas are not necessarily equivalent; in fact, (1) implies (2) under mild assumptions but not *vice versa*. So in our abstract theorems we prefer (1). Indeed, even if (2) implies (1) in all reasonable instances, why postpone for the instantiation time any fact that we can show abstractly?

Applying the Economy principle not only stocks up generality for instantiations, but also accurately outlines trade-offs: How much does it cost (in terms of other added assumptions) to improve the conclusion, or to weaken an assumption of a theorem? For example, an Economy-based proof of Rosser’s variant of \mathcal{IT}_1 reveals how much arithmetic we must factor in for weakening the ω -consistency assumption into consistency.

Related Work Gödel initially gave a proof of \mathcal{IT}_1 and the rough proof idea of \mathcal{IT}_2 [13]. Hilbert and Bernays gave a first detailed proof of \mathcal{IT}_2 [16]. A vast literature was dedicated to the (re)formulation, proof, and analysis of these results [4, 33, 38, 39]. The now canonical line of reasoning goes through the derivability conditions devised by Bernays and Hilbert [16] and simplified by Löb [23]. These conditions have inspired a new branch of modal logic called provability logic [4]. Jeroslow has argued that, unlike previously believed, one condition is redundant when proving \mathcal{IT}_2 [17].

Kreisel [20] and Jeroslow [17] were the first to study abstract conditions on logics under which the incompleteness theorems apply. Buldt [5] surveys the state of the art focusing on \mathcal{IT}_1 . Our abstract approach, based on generic syntax and provability and truth predicates, resembles the style of institution-independent model theory [9, 14] and our previous work on abstract completeness [3] and completeness of ordered resolution [34]. Dimensions of generality that our formalized work does not (yet) explore include quantifier-free logics [17] and arithmetical hierarchy refinements [19]. Our syntax axiomatization is inspired by algebraic theories of the λ -calculi syntax [11, 12, 29].

In the realm of mechanical proofs, the earliest substantial development was due to Sieg [36], who used a prover based on TEM (Theory of Elementary Meta-Mathematics) to formalize parts of the proofs of both \mathcal{IT}_1 and \mathcal{IT}_2 . But the first full proof of \mathcal{IT}_1 was achieved by Shankar [35] in the Boyer-Moore prover, followed by Harrison in HOL Light [15] and O’Connor in Coq [25]. \mathcal{IT}_2 has only been fully proved recently—by Paulson in Isabelle/HOL [26, 27] (who also proved \mathcal{IT}_1). All these mechanizations target theories over a fixed language in classical FOL: that of arithmetic (Harrison and O’Connor) and that of HF sets or a variation of it (Sieg, Shankar and Paulson). These mechanizations are mostly focused on “getting all the work done” in a particular setting (although Harrison targets a more abstract class of theories in the given language). On their way to \mathcal{IT}_1 , Shankar and O’Connor also prove representability of all partial, respectively primitive recursive functions—important standalone results. Also, there has been work on fully automating parts of the proofs of these theorems [1, 6, 32, 37].

By contrast, we explore conditions that enable different formulations for an abstract logic, where aspects such as recursiveness are below our abstraction level. The two approaches are complementary, and they both contribute to formally taming the complex ramifications of the incompleteness theorems. When instantiating our abstract assumptions to recover and upgrade Paulson’s results, we took advantage of Paulson’s substantial work on proving the many low-level lemmas towards the derivability conditions. More should be done at an abstract level to avoid duplicating some of these laborious lemmas when instantiating the theorems to different logics. This will be future work.

2 Abstract Assumptions

Roughly, the incompleteness theorems are considered to hold for logical theories that (1) contain enough arithmetic and (2) are “effective” in that they themselves can be arithmetized. Our goal is to give a general expression of these favorable conditions. To this end, we identify some logic and arithmetic substrata consisting of structure and axioms that express the containment of (various degrees of) arithmetic more abstractly and flexibly than relative interpretations [41]. We also identify abstract notions of encodings and representability that have just what it takes for a working arithmetization.

2.1 The logical substratum

We start with some unspecified sets of variables (Var, ranged over by x, y, z), terms (Term, ranged over by s, t) and formulas (Fmla, ranged over by φ, ψ, χ). We assume that variables are particular terms, $\text{Var} \subseteq \text{Term}$, and that Var is infinite. Free-variables and substitution operators, FVars and $[-/_]$, are assumed for both terms and formulas. We think of $\text{FVars}(t)$ as the (finite) set of free variables of the term t , and similarly for formulas. We call *sentence* any formula with no free variable, and let Sen denote the set of sentences. We think of $s[t/x]$ as the term obtained from s by the (capture-avoiding) substitution of t for the free occurrences of variable x ; we think of $\varphi[t/x]$ as the formula obtained from φ by the substitution of t for the free occurrences of variable x .

In FOL, terms introduce no bindings, so any occurring variable is free. FOL terms fall under our framework, and so do terms with bindings as in λ -calculi and higher-order logic (HOL). To achieve this degree of inclusiveness while also being able to prove interesting results, we work under some well-behavedness assumptions about the free-variables and substitution operators. For example, free-variables distribute over substitution, $\text{FVars}(\varphi[s/x]) = \text{FVars}(\varphi) - \{x\} \cup \text{FVars}(s)$ if $x \in \text{FVars}(\varphi)$, and substitution is compositional, $\varphi[s_1/x_1][s_2/x_2] = \varphi[s_2/x_2][s_1[s_2/x_2]/x_1]$ if $x_1 \neq x_2$ and $x_1 \notin \text{FVars}(s_2)$. Our extended report [30] contains the full list of our generic syntax axioms.

The incompleteness theorems rely heavily on simultaneous substitution, written $\varphi[t_1/x_1, \dots, t_n/x_n]$, whose properties are tricky to formalize—for example, Paulson’s formalization paper dedicates them ample space [27, 6.2]. To address this problem once and for all generically, we define simultaneous substitution from the single-point substitution, $\varphi[t/x]$, and infer its properties from the single-point substitution axioms. For example, we prove that $\text{FVars}(\varphi[s_1/x_1, \dots, s_n/x_n]) = \text{FVars}(\varphi) \cup \bigcup \{\text{FVars}(s_i) - \{x_i\} \mid i \in \{1, \dots, n\} \text{ and } x_i \in \text{FVars}(\varphi)\}$. The technicalities are delicate: To avoid undesired variable replacements, $\varphi[s_1/x_1, \dots, s_n/x_n]$ must be defined as $\varphi[y_1/x_1] \dots [y_n/x_n][s_1/y_1] \dots [s_n/y_n]$ for some fresh y_1, \dots, y_n , the choice of which we must show to be immaterial. This definition’s complexity is reflected in the properties’ proofs. But again, this one-time effort benefits any “customer” logic: In exchange for a well-behaved single-point substitution, it gets back a well-behaved simultaneous substitution.

We let v_1, v_2, \dots be fixed mutually distinct variables. We write Fmla_k for the set of formulas whose free variables are precisely $\{v_1, \dots, v_k\}$, and Fmla_k^\subseteq for the set of formulas whose variables are among $\{v_1, \dots, v_k\}$. Note that $\text{Fmla}_k \subseteq \text{Fmla}_k^\subseteq$ and $\text{Fmla}_0 = \text{Fmla}_0^\subseteq = \text{Sen}$. Given $\varphi \in \text{Fmla}_k^\subseteq$, we write $\varphi(t_1, \dots, t_n)$ instead of $\varphi[t_1/v_1, \dots, t_n/v_n]$.

In addition to free variables and substitution, our theorems will require formulas to be equipped with term equality (\equiv), Boolean connectives (\perp , \top , \rightarrow , \neg , \wedge , \vee), universal and existential quantifiers (\forall , \exists). In our formalization, we assume a minimal list of the above with respect to intuitionistic logic, and define the rest from this minimal list. They are not assumed to be constructors (syntax builders), but operators on terms and formulas, e.g., $\equiv : \text{Term} \rightarrow \text{Term} \rightarrow \text{Fmla}$, $\perp \in \text{Fmla}$, $\forall : \text{Var} \times \text{Fmla} \rightarrow \text{Fmla}$. This caters for logics that do not have them as primitives. For example, HOL defines all connectives and quantifiers from λ -abstraction and either equality or implication.

We fix a unary relation $\vdash \subseteq \text{Fmla}$ on formulas, called *provability*. We write $\vdash \varphi$ instead of $\varphi \in \vdash$, and say the formula φ is *provable*. Whenever certain formula connectives or quantifiers are assumed present, we will assume that \vdash behaves intuitionistically w.r.t. them—namely, we assume the usual (Hilbert-style) intuitionistic FOL axioms with respect to the abstract connectives and quantifiers. Stronger systems, such as those of classical logic, also satisfy these assumptions.

Consistency, denoted Con , is defined as the impossibility to prove false, namely $\not\vdash \perp$. Another central concept is ω -consistency—we carefully choose a formulation that works intuitionistically, with conclusion reminiscent of Gödel’s negative translation [8]:

OCon : For all $\varphi \in \text{Fmla}_1^\subseteq$, if $\vdash \neg \varphi(n)$ for all $n \in \text{Num}$ then $\not\vdash \neg \neg (\exists x. \varphi(x))$.

Assuming classic deduction in \vdash , this is equivalent to the standard formulation: For all $\varphi \in \text{Fmla}_1^\subseteq$, it is not the case that $\vdash \varphi(n)$ for all $n \in \text{Num}$ and $\vdash \neg (\forall x. \varphi(x))$.

Occasionally, we will consider not only provability but also explicit proofs. We fix a set Proof of (entities we call) *proofs*, ranged over by p, q , and a binary relation between proofs p and sentences φ , written $p \Vdash \varphi$ and read “ p is a proof of φ .” We assume \vdash and \Vdash to be related as expected, in that provability is the same as the existence of a proof:

Rel_\vdash^\Vdash : For all $\varphi \in \text{Sen}$, $\vdash \varphi$ iff there exists $p \in \text{Proof}$ such that $p \Vdash \varphi$.

2.2 The arithmetic substratum

We extend the generic syntax assumptions with a subset $\text{Num} \subseteq \text{Term}$, of *numerals*, ranged over by m, n , which are assumed to be closed, i.e., have no free variables.

Convention 1. In all the shown results we implicitly assume: (1) the generic syntax (free variable and substitution) axioms, (2) at least \rightarrow and \perp plus whatever connectives and quantifiers appear in the statement, (3) closedness of \vdash under intuitionistic deduction rules, and (4) the existence of numerals. Other assumptions (e.g., order-like relation axioms, consistency, standard models, etc.) will be indicated explicitly.

On one occasion, we will assume an order-like binary relation modeled by a formula $\prec \in \text{Fmla}_2$. We write $t_1 \prec t_2$ instead of $\prec(t_1, t_2)$ and $\forall x \prec n. \varphi$ instead of $\forall x. x \prec n \rightarrow \varphi$. It turns out that at our level of abstraction it does not matter whether \prec is a strict or a non-strict order. Indeed, we only require the following two properties, where $x \in M$ denotes $\bigvee_{m \in M} x \equiv m$ and \bigvee expresses the disjunction of a finite set of formulas:

Ord_1 : For all $\varphi \in \text{Fmla}_1$ and $n \in \text{Num}$, if $\vdash \varphi(m)$ for all $m \in \text{Num}$, then $\vdash \forall x \prec n. \varphi(x)$.

Ord_2 : For all $n \in \text{Num}$, there exists a finite set $M \subseteq \text{Num}$ such that $\vdash \forall x. x \in M \vee n \prec x$.

Ord_1 states that if a property φ is provable for all numerals, then its universal quantification bounded by any given numeral n is also provable. Having in mind the arith-

metic interpretation of numerals, it would also make sense to assume a stronger version of Ord_1 , replacing “if $\vdash \varphi(m)$ for all $m \in \text{Num}$ ” by the weaker hypothesis “if $\vdash \varphi(m)$ for all $m \in \text{Num}$ such that $\vdash m < n$ ”. But this stronger version will not be needed.

Ord_2 states that, for any numeral n , any element x in the domain of discourse is either greater than n or equal to one of a finite set M of numerals. If we instantiate our syntax to that of first-order arithmetic, then the natural number model satisfies Ord_1 and Ord_2 when interpreting $<$ as either $<$ or \leq . Moreover, these properties are provable in intuitionistic Robinson arithmetic, again for both $<$ and \leq .

2.3 Encodings and representability

Central in the incompleteness theorems are functions that encode formulas and proofs as numerals, $\langle _ \rangle : \text{Fmla} \rightarrow \text{Num}$ and $\langle _ \rangle : \text{Proof} \rightarrow \text{Num}$. For our abstract results, the encodings are not required to be injective or surjective.

Let A_1, \dots, A_m be sets, and let, for each of them, $\langle _ \rangle : A_i \rightarrow \text{Num}$ be an “encoding” function to numerals. Then, an m -ary relation $R \subseteq A_1 \times \dots \times A_m$ is said to be *represented* by a formula $\textcircled{R} \in \text{Fmla}_m$ if the following hold for all $(a_1, \dots, a_m) \in A_1 \times \dots \times A_m$:

- $(a_1, \dots, a_m) \in R$ implies $\vdash \textcircled{R}(\langle a_1 \rangle, \dots, \langle a_m \rangle)$
- $(a_1, \dots, a_m) \notin R$ implies $\vdash \neg \textcircled{R}(\langle a_1 \rangle, \dots, \langle a_m \rangle)$

Let A be another set with $\langle _ \rangle : A \rightarrow \text{Num}$. An m -ary function $f : A_1 \times \dots \times A_m \rightarrow A$ is said to be *represented* by a formula $\textcircled{f} \in \text{Fmla}_{m+1}$ if for all $(a_1, \dots, a_m) \in A_1 \times \dots \times A_m$:

- $\vdash \textcircled{f}(\langle a_1 \rangle, \dots, \langle a_m \rangle, \langle f(a_1, \dots, a_m) \rangle)$
- $\vdash \forall x, y. \textcircled{f}(\langle a_1 \rangle, \dots, \langle a_m \rangle, x) \wedge \textcircled{f}(\langle a_1 \rangle, \dots, \langle a_m \rangle, y) \rightarrow x \equiv y$

The notion of a function being represented is stronger than that of its graph being represented (as a relation)—but with enough deductive power they are equivalent [38, §16]. We will need an even stronger notion: A function f as above is *term-represented* by an operator $\textcircled{f} : \text{Term}^m \rightarrow \text{Term}$ if $\vdash \textcircled{f}(\langle a_1 \rangle, \dots, \langle a_m \rangle) \equiv \langle f(a_1, \dots, a_m) \rangle$ for all $(a_1, \dots, a_m) \in A_1 \times \dots \times A_m$. When the formula by which a relation/function P is represented or term-represented is irrelevant, we call P *representable* or *term-representable*.

We will also need an enhancement of relation representability: Given $i < m$, we call the representation of an m -ary relation R by \textcircled{R} *i-clean* if $\vdash \neg \textcircled{R}(n_1, \dots, n_m)$ for all numbers n_1, \dots, n_m such that n_i (the i ’th number among them) is outside the image of $\langle _ \rangle$ (i.e., there is no $a \in A_i$ with $n_i = \langle a \rangle$). Cleaness would be trivially satisfied if the encodings were surjective. However, surjectivity is not a reasonable assumption. For example, most of the numeric encodings used in the literature are injective but not surjective.

We let $S : \text{Fmla}_1 \rightarrow \text{Sen}$ be the *self-substitution* function, which sends any $\varphi \in \text{Fmla}_1$ to $\varphi(\langle \varphi \rangle)$, i.e., to the sentence obtained from φ by substituting the encoding of φ for the unique variable of φ . An alternative is the following “soft” version of S , which sends any $\varphi \in \text{Fmla}_1$ to $\exists v_1. v_1 \equiv \langle \varphi \rangle \wedge \varphi$, where v_1 is the single free variable of φ . The soft version yields provably equivalent formulas and has the advantage that it is easier to represent inside the logic, since it does not require formalizing the complexities of capture-avoiding substitution. All our results involving S have been proved for both versions.

We will consider the properties Repr_\neg , Repr_S , and Repr_\perp , stating the representability of the functions \neg and S , and of the relation \Vdash . In addition, Clean_\perp will state that the considered representation of \Vdash is 1-clean, i.e., it is clean on the proof component. For

the representing formulas for the above relations and functions we will use their circled names, \ominus , \oplus , etc.; for example, Repr_{\vdash} means that (1) $p \Vdash \varphi$ implies $\vdash \oplus(\langle p \rangle, \langle \varphi \rangle)$ and (2) $p \nVdash \varphi$ implies $\vdash \neg \oplus(\langle p \rangle, \langle \varphi \rangle)$ for all $p \in \text{Proof}$ and $\varphi \in \text{Sen}$.

2.4 Derivability conditions

Most of our assumptions refer to representability. An important exception is the provability relation \vdash , for which only a weakening of representability is reasonable. Let $\oplus \in \text{Fmla}_1$ be the formula for this task. We consider the following assumptions about \oplus , known as the Hilbert-Bernays-Löb derivability conditions:

- HBL₁: $\vdash \varphi$ implies $\vdash \oplus(\varphi)$ for all $\varphi \in \text{Sen}$.
- HBL₂: $\vdash \oplus(\varphi) \wedge \oplus(\varphi \rightarrow \psi) \rightarrow \oplus(\psi)$ for all $\varphi, \psi \in \text{Sen}$.
- HBL₃: $\vdash \oplus(\varphi) \rightarrow \oplus(\oplus(\varphi))$ for all $\varphi \in \text{Sen}$.

Above and elsewhere, to lighten notation we omit parentheses when instantiating one-variable formulas with encodings of formulas—e.g., writing $\oplus(\varphi)$ instead of $\oplus(\langle \varphi \rangle)$.

HBL₁ states that, if a sentence is provable, then its encoding is also provable inside the representation. HBL₃ is roughly a formulation of HBL₁ “one level up,” inside the proof system \vdash . Finally, note that the provability relation is closed under *modus ponens*, in that $\vdash \varphi$ and $\vdash \varphi \rightarrow \psi$ implies $\vdash \psi$ for all $\varphi, \psi \in \text{Sen}$. Thus, HBL₂ roughly states the same property inside the proof system. In short, the derivability conditions state that the representation of provability acts partly similarly to the provability relation. Note that the representability of “proof of” implies HBL₁, taking $\oplus(x)$ to be $\exists y. \oplus(y, x)$.

Convention 2. We focus on the standard provability representation in this paper: Whenever we assume explicit proofs and representability of “proof of,” the formula \oplus will be defined from \oplus as shown above.

We will also be interested in the following variations of the derivability conditions:

- HBL₄: $\vdash \oplus(\varphi) \wedge \oplus(\psi) \rightarrow \oplus(\varphi \wedge \psi)$ for all $\varphi, \psi \in \text{Sen}$.
- HBL₁[≠]: $\vdash \oplus(\varphi)$ implies $\vdash \varphi$ for all $\varphi \in \text{Sen}$.
- SHBL₃: $\vdash \oplus(t) \rightarrow \oplus(\oplus(t))$ for all closed terms t .
- WHBL₂: $\vdash \varphi \rightarrow \psi$ implies $\vdash \oplus(\varphi) \rightarrow \oplus(\psi)$ for all $\varphi, \psi \in \text{Sen}$.

HBL₄ has a similar flavor as HBL₂, but refers to conjunction: It states that the conjunction introduction rule holds inside the proof system. HBL₁[≠] is the converse of HBL₁. Finally, SHBL₃ is a strengthening of HBL₃ holding for all closed terms and not only those that encode sentences, and (if we assume HBL₁) WHBL₂ is a weakening of HBL₂.

2.5 Standard models

We fix a unary relation $\models \subseteq \text{Sen}$, representing *truth of a sentence in the standard model*. We write $\models \varphi$ instead of $\varphi \in \models$, and read it as “ φ is true.” We consider the assumptions:

- Syn _{\models} : Syntactic entities (logical connectives and quantifiers) handle truth as expected:
 - (1) $\not\models \perp$; (2) for all $\varphi, \psi \in \text{Sen}$, $\models \varphi$ and $\models \varphi \rightarrow \psi$ imply $\models \psi$;
 - (3) for all $\varphi \in \text{Fmla}_1$, if $\models \varphi(n)$ for all $n \in \text{Num}$ then $\models \forall x. \varphi(x)$;
 - (4) for all $\varphi \in \text{Fmla}_1$, if $\models \exists x. \varphi(x)$ then $\models \varphi(n)$ for some $n \in \text{Num}$;
 - (5) for all $\varphi \in \text{Sen}$, $\models \varphi$ or $\models \neg \varphi$.

Soundness (of provability with respect to truth): $\vdash \varphi$ implies $\models \varphi$ for all $\varphi \in \text{Sen}$.

$\text{Syn}_{\models}(1-4)$ only contains a partial description of the syntactic entities' behavior—corresponding to elimination rules for \perp , \rightarrow and \exists and introduction rule for \forall . For our results this suffices. $\text{Syn}_{\models}(5)$ states that standard models decide every sentence.

On his way to formalizing \mathcal{IT}_2 for extensions of the HF set theory, after proving HBL_1 Paulson notes [27, p.21]: “The reverse implication [namely $\text{HBL}_1^{\Leftarrow}$], despite its usefulness, is not always proved.” In his abstract account, Buldt also assumes $\text{HBL}_1^{\Leftarrow}$ in his most general formulation of \mathcal{IT}_1 [5, Theorem 3.1]; that formulation has in mind not necessarily the standard provability representation (our Convention 2), but any formula that weakly represents \vdash , which is acceptable for \mathcal{IT}_1 but not for \mathcal{IT}_2 [2].

We avoid such an \mathcal{IT}_1 versus \mathcal{IT}_2 divergence by remaining focused on the standard provability representation. In this case, for arithmetics and related theories, $\text{HBL}_1^{\Leftarrow}$ cannot be inferred without assuming soundness in the standard model (which Paulson does), or at least ω -consistency. We can depict the situation abstractly, without knowing what standard models look like:

Lemma 3. (1) Assume $\text{Rel}_{\vdash}^{\vdash}$, Repr_{\vdash} , Clean_{\vdash} and OCon . Then $\text{HBL}_1^{\Leftarrow}$ holds.
(2) Assume Soundness and $\text{Syn}_{\models}(1,2,3)$. Then OCon holds.
(3) Assume $\text{Rel}_{\vdash}^{\vdash}$, Repr_{\vdash} , Clean_{\vdash} , Soundness and $\text{Syn}_{\models}(1,2,4)$. Then $\models \oplus\langle\varphi\rangle$ implies $\vdash \varphi$ for all $\varphi \in \text{Sen}$. In particular, $\text{HBL}_1^{\Leftarrow}$ holds.

Thus, staying in a proof-theoretic world, ω -consistency ensures $\text{HBL}_1^{\Leftarrow}$ if the “proof of” relation is cleanly represented (1). In turn, ω -consistency is ensured by minimal semantic requirements, including the soundness of provability (2). Finally, putting together representability and semantics, we can infer something stronger than $\text{HBL}_1^{\Leftarrow}$: That the mere truth (and not just the provability) of a sentence's provability representation implies the provability of the sentence itself (3).

It follows from either points (1,2) or point (3) of the lemma that, in the presence of standard models and soundness, clean representability of the “proof of” relation implies $\text{HBL}_1^{\Leftarrow}$; and recall that it also implies HBL_1 . So it implies an “iff” version of HBL_1 : $\vdash \varphi$ if and only if $\vdash \oplus\langle\varphi\rangle$. Interestingly, a converse of this implication also holds. To state it, we initially assume there is no “outer” notion of proof (i.e., no set Proof and no relation \Vdash), but only an “inner” one, given by a formula $P \in \text{Fmla}_2$ such that:

$\text{Rel}_{\oplus}^P: \vdash \oplus\langle\varphi\rangle \leftrightarrow \exists x. P(x, \langle\varphi\rangle)$.
 $\text{Compl}_P: \models P(n, \langle\varphi\rangle)$ implies $\vdash P(n, \langle\varphi\rangle)$ for all $n \in \text{Num}$ and $\varphi \in \text{Sen}$.
 $\text{Compl}_{\neg P}: \models \neg P(n, \langle\varphi\rangle)$ implies $\vdash \neg P(n, \langle\varphi\rangle)$ for all $n \in \text{Num}$ and $\varphi \in \text{Sen}$.

Rel_{\oplus}^P is the inner version of $\text{Rel}_{\vdash}^{\vdash}$: It expresses that, *inside the representation*, proofs and provability are connected as expected. Compl_P and $\text{Compl}_{\neg P}$ state that provability is complete on P statements about formula encodings, as well as their negations; in traditional settings, this is true thanks to P being a bounded arithmetical formula (\mathcal{A}_0). Now the converse result states that, thanks to (standard models and) the “iff” version of HBL_1 , we can define an outer notion of proof that is represented by the inner notion P :

Lemma 4. Assume Rel_{\oplus}^P , Compl_P , $\text{Compl}_{\neg P}$, Soundness, $\text{Syn}_{\models}(4,5)$, HBL_1 and $\text{HBL}_1^{\Leftarrow}$. Take $\text{Proof} = \text{Num}$ and define \Vdash by $n \Vdash \varphi$ iff $\vdash P(n, \langle\varphi\rangle)$. Then $\text{Rel}_{\vdash}^{\vdash}$, Repr_{\vdash} and Clean_{\vdash} hold, with \vdash being represented by P .

3 Abstract Incompleteness Theorems

After last section's preparations, we are now ready to discuss different versions of the incompleteness theorems and their major lemmas, based on alternative assumptions.

3.1 Diagonalization

The formula diagonalization technique (due to Gödel and Carnap [7]) yields “self-referential” sentences. All we need for it to work is the representability of substitution.

Prop 5. Assuming Repr_S , for all $\psi \in \text{Fmla}_1$ there exists $\varphi \in \text{Fmla}_1$ with $\vdash \varphi \leftrightarrow \psi(\varphi)$.

A sentence $\varphi \in \text{Sen}$ is called a *Gödel sentence* if $\vdash \varphi \leftrightarrow \neg \bigoplus(\varphi)$; it is called a *Rosser sentence* if $\vdash \varphi \leftrightarrow \neg (\exists x. \bigoplus(x, \langle \varphi \rangle) \wedge \text{RosserTwist}(x, \langle \varphi \rangle))$, where we define $\text{RosserTwist}(x, y) = \forall x'. x' \prec x \rightarrow \forall y'. \bigodot(y, y') \rightarrow \neg \bigoplus(x', y')$. The existence of Gödel and Rosser sentences follows immediately from diagonalization.

Prop 6. Assuming Repr_S , there exist Gödel and Rosser sentences.

Thus, any Gödel sentence is provably equivalent to the negation of its own provability; in Gödel's words [13], it “says about itself that it is not provable.” A Rosser sentence φ asserts its own unprovability in a weaker fashion: Rather than saying “Myself, φ , am not provable” (i.e., “it is not the case that there exists a proof p of φ ”), it says “it is not the case that there exists a proof p of φ such that, for all smaller proofs q , q is not a proof of $\neg \varphi$.” Here, “smaller” refers to the order the encoding of proofs as numerals imposes.

3.2 The incompleteness theorems

\mathcal{IT}_1 identifies sentences that are neither provable nor disprovable—which often holds for Gödel and Rosser sentences with the help of a provability relation satisfying HBL_1 .

Prop 7. Assume Con and HBL_1 . Then $\not\vdash G$ for all Gödel sentences G .

For showing that the Gödel sentences are not disprovable, a standard route is to assume explicit proofs, strengthen the consistency assumption to ω -consistency, and strengthen HBL_1 to representability of the “proof of” relation.

Prop 8. Assume OCon , $\text{Rel}_{\vdash}^{\vdash}$, Repr_{\vdash} , Clean_{\vdash} . Then $\not\vdash \neg G$ for all Gödel sentences G .

Proof. Let G be a Gödel sentence. We prove $\not\vdash \neg G$ by contradiction. Assume (1) $\vdash \neg G$.

- By consistency (which is implied by OCon), we obtain $\not\vdash G$.
- From this and $\text{Rel}_{\vdash}^{\vdash}$, we obtain $p \not\vdash G$ for all $p \in \text{Proof}$.
- From this, Repr_{\vdash} and Clean_{\vdash} , we obtain $\vdash \neg \bigoplus(n, \langle G \rangle)$ for all $n \in \text{Num}$.
- From this and OCon , we obtain $\not\vdash \neg \neg \exists x. \bigoplus(x, \langle G \rangle)$, i.e., $\not\vdash \neg \neg \bigoplus(\langle G \rangle)$.
- Hence, since G is a Gödel sentence, we obtain $\not\vdash \neg G$, which contradicts (1). \square

While the line of reasoning in the above proof is mostly well-known, it contains two subtle points about which the literature is not explicit (due to the usual focus on classical first-order arithmetic and particular choices of encodings).

First, we must assume the representation of the “proof of” relation to be 1-*clean*, i.e., clean with respect to the proof component. Indeed, the argument crucially relies on converting the statement “ $p \not\vdash G$ for all $p \in \text{Proof}$ ” into “ $\vdash \neg \bigoplus(n, \langle G \rangle)$ for all $n \in \text{Num}$,” which is only possible for 1-clean encodings. This assumption will be repeatedly needed

in later results. By contrast, cleanness is never required with respect to the sentence component of “proof of” or for the provability relation (which only involves sentence encodings). In short, cleanness is only needed for proofs, not for sentences.

Second, to reach the desired contradiction for our intuitionistic proof system \vdash , from “ $\vdash \neg \oplus(n, \langle G \rangle)$ for all $n \in \text{Num}$ ” it is not sufficient to employ standard ω -consistency, which would only give us $\not\vdash \exists x. \oplus(x, \langle G \rangle)$, i.e., $\not\vdash \oplus \langle G \rangle$; the last together with $\vdash G \leftrightarrow \neg \oplus \langle G \rangle$ would be insufficient for obtaining $\not\vdash \neg G$. However, our stronger version of ω -consistency, OCon, does the trick. \mathcal{IT}_1 now follows by putting together Props. 6–8:

Theorem 9. (\mathcal{IT}_1) Assume OCon, $\text{Rel}_{\vdash}^{\vdash}$, Repr_{\vdash} , Clean_{\vdash} , and Repr_S . Then:

- (1) There exists a Gödel sentence. (2) $\not\vdash G$ and $\not\vdash \neg G$ for all Gödel sentences G .

Rosser’s contribution to \mathcal{IT}_1 was an ingenious trick for weakening the ω -consistency assumption into plain consistency—as such, it is usually seen as a *strict improvement* over Gödel’s version. While this is true for the concrete case of FOL theories extending arithmetic, from an abstract perspective the situation is more nuanced: The improvement is achieved at the cost of asking more from the logic. Our framework makes this trade-off clearly visible. The idea is to use Rosser sentences instead of Gödel sentences to “repair” the ω -consistency assumption of Theorem 9 (inherited from Prop. 8):

Theorem 10. (\mathcal{IT}_1 à la Rosser) Assume Con, Ord₁, Ord₂, Repr_¬, $\text{Rel}_{\vdash}^{\vdash}$, Repr_{\vdash} , Clean_{\vdash} , and Repr_S . Then:

- (1) There exists a Rosser sentence. (2) $\not\vdash R$ and $\not\vdash \neg R$ for all Rosser sentences R .

Highlighted is the assumption trade-off between the two versions: Rosser’s weakening of ω -consistency into consistency is paid by additionally assuming representability of negation and an order-like relation satisfying Ord₁ and Ord₂. Certainly, negation representability is not a big price, since for concrete logics this tends to be a lemma that is anyway needed when proving HBL₁. On the other hand, the ordering assumptions seem to be a significant generality gap in favor of Gödel’s version. A clear manifestation of this gap is in our inference of a semantic version of \mathcal{IT}_1 —which we obtain from Theorem 9 with the help of Lemmas 3(2) and 4:

Theorem 11. (Semantic \mathcal{IT}_1) Assume $\text{Rel}_{\oplus}^{\oplus}$, Compl_P, Compl_{¬P}, Soundness, Syn_{\vdash} , HBL₁, $\text{HBL}_1^{\Leftarrow}$, and Repr_S . Then:

- (1) There exists a Gödel sentence. (2) $\models G$, $\not\vdash G$, and $\not\vdash \neg G$ for all Gödel sentences G .

We have highlighted the assumptions specific to the semantic treatment. They replace OCon, $\text{Rel}_{\vdash}^{\vdash}$, Repr_{\vdash} and Clean_{\vdash} from the proof-theoretic Theorem 9. Also highlighted is the additional fact concluded: that the Gödel sentences are true.

We have inferred the semantic version from Gödel’s proof-theoretic version (Theorem 9), and not from Rosser’s variation (Theorem 10). This is because in the semantic version ω -consistency comes for free (from Lemma 3(2)). By contrast, for deploying Rosser’s version we would need to explicitly consider the order-like relation with its own hypotheses. This would have led to a *strictly less general* abstract result (if we ignore the difference in the way Gödel and Rosser sentences are actually defined).

The semantic \mathcal{IT}_1 relies on $\text{HBL}_1^{\Leftarrow}$. If we commit to classical logic (i.e., assume

$\vdash \neg \neg \varphi \rightarrow \varphi$), we can more directly show, taking advantage of $\text{HBL}_1^{\Leftarrow}$, that the Gödel sentences are not disprovable, which immediately proves \mathcal{IT}_1 :

Theorem 12. (Classical \mathcal{IT}_1) Assume classical logic, Con, HBL_1 , $\text{HBL}_1^{\Leftarrow}$, Repr_S . Then:

- (1) There exists a Gödel sentence. (2) $\nvdash G$ and $\nvdash \neg G$ for all Gödel sentences G .

Classical logic also offers two alternatives to our semantic Theorem 11 (where the second is strictly more general than the first):

Theorem 13. (Classical Semantic \mathcal{IT}_1) The conclusions of Theorem 11 still hold if we assume classical logic and perform either of the following changes in its assumptions: (1) remove $\text{Compl}_{\neg P}$, or (2) replace Rel_{\oplus}^P , Compl_P and $\text{Compl}_{\neg P}$ with “ $\models \oplus(\varphi)$ implies $\vdash \varphi$ for all $\varphi \in \text{Sen}$.”

Even though \mathcal{IT}_1 needs a predicate \oplus that satisfies HBL_1 (and sometimes also $\text{HBL}_1^{\Leftarrow}$, meaning that it weakly represents provability), its conclusion, the existence of undecided sentences, is meaningful regardless of whether \oplus *adequately expresses provability*. By contrast, the meaning of \mathcal{IT}_2 ’s conclusion, the theory cannot prove its own consistency, relies on this (non-mathematical) “intensional” assumption [2]. In this case, consistency is adequately expressed by the sentence $\neg \oplus(\perp)$. The standard formulation (and proof) of \mathcal{IT}_2 uses all three derivability conditions:

Theorem 14. (\mathcal{IT}_2) Assume Con, HBL_1 , HBL_2 , HBL_3 and Repr_S . Then $\nvdash \neg \oplus(\perp)$.

3.3 Jeroslow’s approach

Next we study an alternative line of reasoning due to Jeroslow [17], often cited as a simplification of the canonical route to prove \mathcal{IT}_2 [33, 38, 39]. To study its features and pitfalls, we need some standard notation used by Jeroslow. A *pseudo-term* is a formula $\varphi \in \text{Fmla}_{m+1}$ expressing a provably functional relation via “exists unique”: $\vdash \forall x_1, \dots, x_m. \exists! y. \varphi(x_1, \dots, x_m, y)$. We only discuss the case $m = 2$; the general case is similar.

Notation 15. Given a pseudo-term $\varphi \in \text{Fmla}_2$, we treat it as if it is a one-variable term:

- for any terms s and t , we write $t \equiv \varphi(s)$ instead of $\varphi(s, t)$;
- for any term s and formula $\psi \in \text{Fmla}_1$, we write $\psi(\varphi(s))$ instead of $\exists y. \varphi(s, y) \wedge \psi(y)$.

This notation smoothly integrates pseudo-terms with terms: If $\vdash t \equiv \varphi(s)$ and $\vdash \psi(\varphi(s))$ then $\vdash \psi(t)$, where $\psi(t)$ denotes actual substitution of terms in formulas.

Jeroslow relies on an abstract class of m -ary functions, $\mathcal{F}_m \subseteq \text{Num}^m \rightarrow \text{Num}$, for all arities $m \in \mathbb{N}$, on which he considers the following assumptions:

$\text{Repr}_{\mathcal{F}}$: Every $f \in \mathcal{F}_m$ is represented by some pseudo-term $\textcircled{f} \in \text{Fmla}_{m+1}$ under the identity encoding $\text{Num} \rightarrow \text{Num}$.

CapN : Some $N \in \mathcal{F}_1$ correctly captures negation: $N(\varphi) = \langle \neg \varphi \rangle$ for all $\varphi \in \text{Sen}$.

CapSS : Some $\text{ssap} : \text{Fmla}_1 \rightarrow \mathcal{F}_1$ correctly captures substituted self-application:

$$\text{ssap } \psi \langle \textcircled{f} \rangle = \langle \psi(\textcircled{f}(\textcircled{f})) \rangle \text{ for all } \psi \in \text{Fmla}_1 \text{ and } f \in \mathcal{F}_1.$$

In CapSS , following Jeroslow we employed Notation 15 taking advantage of the fact that \textcircled{f} are pseudo-terms: The highlighted text denotes $\exists y. \textcircled{f}(\langle \textcircled{f} \rangle, y) \wedge \psi(y)$. Moreover, using the same notation, the statement of $\text{Repr}_{\mathcal{F}}$ for some $f \in \mathcal{F}_1$ and $n \in \text{Num}$ would be written as $\vdash f(n) \equiv \textcircled{f}(n)$. Similarly, combining CapN with the instance of $\text{Repr}_{\mathcal{F}}$, we obtain a fact that can be written as $\vdash \langle \neg \varphi \rangle \equiv \textcircled{N}(\varphi)$.

When our logical theory is a recursive extension of Robinson arithmetic and $\text{Num} = \mathbb{N}$, \mathcal{F}_m could be the set of m -ary computable functions. Then every $f \in \mathcal{F}_m$ would indeed be represented by a formula $\langle f \rangle$. Moreover, assuming a computable and injective encoding of formulas, $\langle _ \rangle : \text{Fmla}_1 \rightarrow \mathbb{N}$, we can take $N : \mathbb{N} \rightarrow \mathbb{N}$ to be the following computable function: Given input n , it checks if n has the form $\langle \varphi \rangle$; if so, it returns $\langle \neg \varphi \rangle$; if not, it returns any value (e.g., 0). And $\text{ssap } \psi$ can be defined similarly, obtaining the desired property for every $\varphi \in \text{Fmla}_2$, not necessarily of the form $\langle f \rangle$. In short, Jeroslow's assumptions cover arithmetic (but also potentially many other systems).

At the heart of Jeroslow's approach lies an alternative diagonalization technique, producing *term* fixpoints, not just formula fixpoints:

Lemma 16. Assume CapSS and $\text{Repr}_{\mathcal{F}}$ and let $\psi \in \text{Fmla}_1$. Then there exists a closed pseudo-term t such that $\vdash t \equiv \langle \psi(t) \rangle$. Moreover, taking $\varphi = \psi(t)$, we have $\vdash \varphi \leftrightarrow \psi\langle \varphi \rangle$.

Proof. Let $f = \text{ssap } \psi$ and $t = \langle f \rangle \langle f \rangle$. From CapSS , we obtain $f\langle f \rangle = \langle \psi(\langle f \rangle \langle f \rangle) \rangle$. From this and $\text{Repr}_{\mathcal{F}}$, we obtain $\vdash \langle f \rangle \langle f \rangle \equiv \langle \psi(\langle f \rangle \langle f \rangle) \rangle$, i.e., $\vdash t \equiv \langle \psi(t) \rangle$. With the equality rules, we obtain $\vdash \psi(t) \leftrightarrow \psi(\langle \psi(t) \rangle)$, i.e., $\vdash \varphi \leftrightarrow \psi\langle \varphi \rangle$. \square

This lemma offers us Gödel and Rosser sentences, which can be used like in Sections 3.1 and 3.2, leading to corresponding variants of \mathcal{IT}_1 . But Jeroslow's main innovation affects \mathcal{IT}_2 : While traditionally \mathcal{IT}_2 requires all three derivability conditions, Jeroslow's version does not make use of the second, HBL_2 :

Theorem 17. (\mathcal{IT}_2 à la Jeroslow) Assume Con , HBL_1 , SHBL_3 , $\text{Repr}_{\mathcal{F}}$, CapN , CapSS . Then $\nvdash \text{jcon}$, where jcon denotes $\forall x. \neg (\oplus(x) \wedge \oplus(\mathbb{N}(x)))$.

Like with Rosser's trick, we analyze this innovation's trade-offs from an abstract perspective. A first trade-off is in the employment of a stronger version of the third condition, SHBL_3 (extended to affect all closed pseudo-terms via Notation 15). Another is in the way consistency is expressed in the logic. Jeroslow does not conclude $\nvdash \neg \oplus(\perp)$, but something more elaborate, namely $\nvdash \text{jcon}$. While the formula $\neg \oplus(\perp)$ internalizes the statement $\nvdash \perp$, jcon internalizes the equivalent statement “for all φ , it is not the case that $\vdash \varphi$ and $\vdash \neg \varphi$.” But are the internalizations themselves equivalent, i.e., is it the case that $\vdash \neg \oplus(\perp)$ iff $\vdash \text{jcon}$? This surely holds for many concrete logics, but it is one direction that we can infer logic-independently: Assuming HBL_1 , $\text{Repr}_{\mathcal{F}}$ and CapN , $\vdash \text{jcon}$ implies $\vdash \neg \oplus(\perp)$. And it seems we cannot infer the other direction without knowing what \oplus looks like more concretely. Therefore, $\nvdash \neg \oplus(\perp)$, the conclusion of the original \mathcal{IT}_2 , is *abstractly stronger than*, hence *preferable to* $\nvdash \text{jcon}$. In short, Jeroslow somewhat weakens the theorem's conclusion.

Let us now look at (a slight rephrasing of) Jeroslow's proof:

Proof of Theorem 17. We assume (1) $\vdash \text{jcon}$ and aim to reach a contradiction.

- Applying Lemma 16 to $\oplus(\mathbb{N}(x))$, obtain a closed term t where (2) $\vdash t \equiv \langle \oplus(\mathbb{N}(t)) \rangle$.
- By SHBL_3 applied to $\mathbb{N}(t)$, we obtain $\vdash \oplus(\mathbb{N}(t)) \rightarrow \oplus(\oplus(\mathbb{N}(t)))$.
- From (2) and the equality rules, we obtain $\vdash \oplus(\mathbb{N}(t)) \rightarrow \oplus(\mathbb{N}(\oplus(\mathbb{N}(t))))$.
- The last two facts give us $\vdash \varphi \rightarrow \oplus(\varphi) \wedge \oplus(\mathbb{N}(\varphi))$, where φ denotes $\oplus(\mathbb{N}(t))$.
- On the other hand, (1) instantiated with $\langle \varphi \rangle$ gives us $\vdash \neg (\oplus(\varphi) \wedge \oplus(\mathbb{N}(\varphi)))$.
- From the last two facts, we obtain (3) $\vdash \neg \varphi$.

- With HBL_1 , we obtain $\vdash \oplus(\neg\varphi)$ and with CapN and $\text{Repr}_{\mathcal{F}}$, we obtain $\vdash \oplus(\mathbb{N}(\varphi))$.
- From (2) and the equality rules, we obtain $\vdash \oplus(\mathbb{N}(\oplus(\mathbb{N}(t)))) \rightarrow \oplus(\mathbb{N}(t))$, i.e., $\vdash \oplus(\mathbb{N}(\varphi)) \rightarrow \varphi$
- From the last two facts, we obtain $\vdash \varphi$. With (3) this contradicts (1). \square

A first major observation is that, under the stated assumptions, the above proof is *incorrect*. It uses an implicit assumption, hidden under Notation 15: When we disambiguate the notation, we see that Lemma 16 gives us a pseudo-term t that does not exactly satisfy (1) $\vdash t \equiv \langle \psi(t) \rangle$ (which is what the theorem's proof needs), but something weaker, namely (2) $\vdash t \equiv \langle \chi \rangle$, where χ is $\vdash \exists x. \mathcal{F}(\langle \mathcal{F} \rangle, x) \wedge \psi(x)$. And although $\vdash \chi \leftrightarrow \psi(t)$, we still cannot infer (1) from (2), unless *the encodings of provably equivalent formulas are assumed provably equal*. But this assumption is unreasonable: Usually formula equivalence is undecidable, so no computable encoding can achieve that. (Incidentally, this problem is also the reason why we need SHBL_3 instead of HBL_3 : In the proof's application of SHBL_3 to obtain $\vdash \oplus(\mathbb{N}(t)) \rightarrow \oplus(\oplus(\mathbb{N}(t)))$, we cannot work with $\langle \neg\varphi \rangle$ instead of $\mathbb{N}(t)$, even though $\vdash \langle \neg\varphi \rangle \equiv \mathbb{N}(t)$.)

To repair that, we can replace representation by pseudo-terms with actual term-representation. More precisely (also factoring in the observation that Jeroslow's proof does not need \mathcal{F}_n for all n , but \mathcal{F}_1 suffices), we change $\text{Repr}_{\mathcal{F}}$ into:

$\text{Repr}_{\mathcal{F}}$: Every $f \in \mathcal{F}_1$ is term-represented, under the identity encoding $\text{Num} \rightarrow \text{Num}$, by some \mathcal{F} taken from a set $\text{Ops} \subseteq (\text{Term} \rightarrow \text{Term})$ for which an encoding as numerals $\langle _ \rangle : \text{Ops} \rightarrow \text{Num}$ is given, and such that $\text{FVars}(g(t)) = \text{FVars}(t)$ and $(g(t))[s/x] = g(t[s/x])$ for all $g \in \text{Ops}$, $s, t \in \text{Term}$ and $x \in \text{Var}$.

(In concrete logics, the elements of Ops can be constructors or derived operators on terms.) Then CapSS , Lemma 16, and all proofs work with terms rather than pseudo-terms and everything becomes formally correct. In summary, Jeroslow's approach to \mathcal{IT}_2 seems to fail for pseudo-terms representing computable functions, but to require actual terms. This usually means that the logic has built-in Skolem symbols and axioms.

Finally, let us see what it takes to alleviate the second trade-off: from $\not\vdash \text{jcon}$ to the more desirable $\not\vdash \neg \oplus(\perp)$. We see that Theorem 17's proof uses $\vdash \text{jcon}$ not at jcon 's full generality but only instantiated with formula encodings, which thanks to $\text{Repr}_{\mathcal{F}}$ and CapN would follow from $(*) \vdash \neg (\oplus(\varphi) \wedge \oplus(\neg\varphi))$. And it only takes WHBL_2 (a weaker version of HBL_2) and HBL_4 to prove $\vdash (\oplus(\varphi) \wedge \oplus(\neg\varphi)) \rightarrow \oplus(\perp)$, allowing us to infer $(*)$ from $\vdash \neg \oplus(\perp)$; meaning that the latter could have been used. We obtain:

Theorem 18. If in the (corrected) Theorem 17 we additionally assume WHBL_2 and HBL_4 , its conclusion can be upgraded to $\not\vdash \neg \oplus(\perp)$.

Whether WHBL_2 and HBL_4 are a good trade-off for HBL_2 will of course depend on the logic's specificity, in particular, on its primitive rules of inference.

Jeroslow presented his approach for an abstract logical theory over a FOL language, which is not necessarily a FOL theory—so it found a natural fit in our generic framework. To our knowledge, very few subsequent authors present Jeroslow's approach rigorously, and none at its original level of generality. Smith's monograph gives a rigorous account for arithmetic [38, §33], silently performing the correction we have shown here, but failing to detect the need for SHBL_3 instead of HBL_3 (which Jeroslow had noticed). A mechanical prover is of invaluable help with detecting such nuances and pitfalls.

Summary Using our generic infrastructure (Section 2), we have formally proved several abstract incompleteness results. They include four versions of \mathcal{IT}_1 :

- Gödel’s original \mathcal{IT}_1 (Theorem 9) and an \mathcal{IT}_1 based on classical logic (Theorem 12) required the formalization of some well-known arguments without change.
- Rosser’s \mathcal{IT}_1 (Theorem 10) involved the generalization of a well-known argument: distilling two abstract conditions, Ord_1 and Ord_2 .
- Novel semantic variants of \mathcal{IT}_1 (Theorems 11 and 13) were born from abstractly connecting standard models, HBL_1 ’s “iff” version, and proof representability.

They also include two versions of \mathcal{IT}_2 :

- The standard \mathcal{IT}_2 based on the three derivability conditions (Theorem 14) again only required formalizing a well-known argument.
- The alternative, Jeroslow-style \mathcal{IT}_2 (Theorems 17 and 18) involved a detailed analysis and correction of an existing abstract result.

4 Instances of the Abstract Results

We first validate the assumptions about our abstract logic and arithmetic:

Prop 19. (1) Any FOL theory that extends Robinson arithmetic or the HF set theory satisfies all the axioms in our logical and arithmetical substrata (in Sections 2.1, 2.2).
 (2) If, in addition, the theory is sound, then, together with its corresponding standard model, it also satisfies all our model-theoretic axioms (in Section 2.5).

In particular, point (2) shows that our discussion of standard models applies equally well to \mathbb{N} and the datatype of HF sets. (In the latter case, Num becomes the entire set of closed terms, so that numerals can denote arbitrary HF sets. This shows the versatility of our abstract concept of numeral.) Then we instantiate three of our main theorems:

Theorem 20. (1) Any FOL theory that extends the HF set theory with a finite set of axioms and `is sound in the standard HF set model` satisfies the hypotheses of Theorems 13 and 14. Hence \mathcal{IT}_1 (semantic version) and \mathcal{IT}_2 hold for it.
 (2) Any FOL theory that extends the HF set theory with a finite set of axioms and `is consistent` satisfies the hypotheses of Theorem 14. Hence \mathcal{IT}_2 holds for it.

These instances are heavily based on the lemmas proved by Paulson in his formalization of \mathcal{IT}_1 and \mathcal{IT}_2 [26, 27], who follows and corrects Świerczkowski’s detailed informal account [40]. Point (1) is a restatement of Paulson’s formalized results: theorems *Goedel_I* and *Goedel_II* in [27]. (His theorems also assume consistency, but that is redundant: Consistency follows from his underlying soundness assumption.)

By contrast, point (2) is an upgrade of Paulson’s *Goedel_II*, applicable to any consistent, though possibly unsound theory. This stronger version is in fact \mathcal{IT}_2 ’s standard form, free from any model-theoretic considerations. Paulson had proved both HBL_1 and $\text{HBL}_1^{\Leftarrow}$ taking advantage of soundness, so we needed to discard $\text{HBL}_1^{\Leftarrow}$ and re-prove HBL_1 by replacing any semantic arguments with proofs within the HF calculus. We also removed all invocations of a convenient “truth implies provability for Σ -sentences” lemma, which depended on soundness due to Paulson’s choice of Σ -sentence definition.

This instantiation process has offered important feedback into the abstract results. A formal development such as ours is (largely) immune to reasoning errors, but not

to missing out on useful pieces of generality. We experienced this firsthand with our assumptions about substitution. An *a priori* natural choice was to assume representability of the numeral substitution $Sb : Fmla_1 \times Num \rightarrow Sen$ (defined as $Sb(\varphi, n) = \varphi(n)$), part of which means $(1) \vdash \mathbb{S}b(\langle \varphi \rangle, n, Sb(\varphi, n))$. But Paulson had instead proved $(2) \vdash \mathbb{S}b(\langle \varphi \rangle, \langle n \rangle, Sb(\varphi, n))$. The key difference from (1) is that (2) applies the term encoding function $\langle _ \rangle : Term \rightarrow Num$ to numerals as well (as particular terms); and since his $\langle _ \rangle$ function is injective, it is far from the case that $\langle n \rangle = n$ for all numerals n . Paulson’s version makes more sense than ours when building the results bottom-up: Representability should not discriminate numerals, but filter them through the encodings like other terms. However, top-down our version also made sense: It yielded the incompleteness theorems under reasonable assumptions, which do hold, by the way, for the HF set theory—even though in a bottom-up development one is unlikely to prove them. We resolved this discrepancy through a common denominator: the representability of self-substitution $S : Fmla_1 \rightarrow Sen$ (Section 2.3), which made our results more general.

Paulson’s formalization has also inspired our abstract treatment of standard models (Section 2.5). Since Paulson proves HBL_1^{\Leftarrow} and uses classical logic, an obvious “port of entry” of his \mathcal{IT}_2 into our framework is Theorem 12. But this theorem tells us nothing about the Gödel sentences’ truth. Delving deeper into Paulson’s proof, we noted that he (unconventionally) completely avoids $Repr_{\Vdash}$, and does not even define \Vdash . This raised the question of whether HBL_1^{\Leftarrow} and $Repr_{\Vdash}$ are somehow interchangeable in the presence of standard models—and we found that they indeed are, under mild assumptions about truth. Incidentally, these assumptions were also sufficient for establishing the Gödel sentences’ truth, leading to our semantic \mathcal{IT}_1 (Theorem 11). However, Theorem 11 was not easy to instantiate to Paulson’s \mathcal{IT}_1 . All its assumptions were easy to prove, except for $Compl_{\neg P}$. Whereas Paulson proved that his proof-of predicate is a Σ -formula (which implies $Compl_P$ by Σ -completeness), he did not prove the same for its negation (which would imply $Compl_{\neg P}$). We are confident that this is true (any reasonable proof-of predicate is a Δ -formula), but we leave the laborious formal proof of this fact as future work. Instead, we recovered Paulson’s result as an instance of our Theorem 13.

As future work, we will consider even more general variants of our semantic Theorems 11 and 13, as in Smorynski’s account [39]: by distinguishing between a sound “base” provability relation \vdash_0 and an extension \vdash only required to be consistent or ω -consistent. For example, \vdash_0 could be deduction in HF set theory or a weaker theory and \vdash deduction in a consistent (not necessarily sound) extension of the HF set theory. This two-layered approach would have also benefited Paulson’s original formalization.

Many other logics and logical theories satisfy our theorems’ assumptions. We do *not* require the logic to be reducible to a single syntactic category of formulas, $Fmla$, a single syntactic judgment, \vdash , etc.; but only that such (well-behaved) formulas, provability relation, etc. are identifiable as part of that logic, e.g., localized to a given type and/or relativised by a given predicate. This allows our framework to capture most variants of higher-order logic and type theory (including the variant underlying Isabelle/HOL itself [21, 22]), and also, we believe, many of the logics surveyed by Buldt [5], including non-classical and fuzzy. But enabling “mass instantiation” that is both formal and painless requires more progress on the agenda we started here: recognizing reusable construction and proof patterns and formalizing them as abstract results.

Acknowledgments. We thank Bernd Buldt for his patient explanations on material in his monograph, and the reviewers for insightful comments and suggestions.

References

1. Ammon, K.: An automatic proof of Gödel’s incompleteness theorem. *Artif. Intell.* 61(2), 291–306 (1993)
2. Auerbach, D.: Intensionality and the Gödel theorems. *Philosophical Studies: An International Journal for Philosophy in the Analytic Tradition* 48(3), 337–351 (1985)
3. Blanchette, J.C., Popescu, A., Traytel, D.: Unified classical logic completeness—A coinductive pearl. In: Demri, S., Kapur, D., Weidenbach, C. (eds.) *IJCAR 2014*. LNCS, vol. 8562, pp. 46–60. Springer (2014)
4. Boolos, G.: *The Logic of Provability*. Cambridge University Press (1993)
5. Buldt, B.: The scope of Gödel’s first incompleteness theorem. *Logica Universalis* 8(3), 499–552 (2014)
6. Bundy, A., Giunchiglia, F., Villaflorida, A., Walsh, T.: An incompleteness theorem via abstraction. Tech. rep., Istituto per la Ricerca Scientifica e Tecnologica, Trento (1996)
7. Carnap, R.: Logische syntax der sprache. *Philosophical Review* 44(4), 394–397 (1935)
8. Davis, M.: *The Undecidable: Basic Papers on Undecidable Propositions, Unsolvability Problems, and Computable Functions*. Dover Publication (1965)
9. Diaconescu, R.: *Institution-independent Model Theory*. Birkhäuser, 1st edn. (2008)
10. Feferman, S., Dawson, Jr., J.W., Kleene, S.C., Moore, G.H., Solovay, R.M., van Heijenoort, J. (eds.): *Kurt Gödel: Collected Works*. Vol. 1: Publications 1929–1936. Oxford University Press (1986)
11. Fiore, M.P., Plotkin, G.D., Turi, D.: Abstract syntax and variable binding. In: *Logic in Computer Science (LICS) 1999*, pp. 193–202. IEEE Computer Society (1999)
12. Gabbay, M.J., Mathijssen, A.: Nominal (universal) algebra: Equational logic with names and binding. *J. Log. Comput.* 19(6), 1455–1508 (2009)
13. Gödel, K.: Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I. *Monatshefte für Mathematik und Physik* 38(1), 173–198 (1931)
14. Goguen, J.A., Burstall, R.M.: Institutions: Abstract model theory for specification and programming. *J. ACM* 39(1), 95–146 (1992)
15. Harrison, J.: HOL Light proof of Gödel’s first incompleteness theorem, located at <http://code.google.com/p/hol-light/>, directory Arithmetic
16. Hilbert, D., Bernays, P.: *Grundlagen der Mathematik*, Vol. II. Springer-Verlag (1939)
17. Jeroslow, R.G.: Redundancies in the Hilbert-Bernays derivability conditions for Gödel’s second incompleteness theorem. *J. Symb. Log.* 38(3), 359–367 (1973)
18. Kaliszyk, C., Urban, J.: HOL(y)Hammer: Online ATP service for HOL light. *Mathematics in Computer Science* 9(1), 5–22 (2015)
19. Kikuchi, M., Kurahashi, T.: Generalizations of Gödel’s incompleteness theorems for Σ_n -definable theories of arithmetic. *Rew. Symb. Logic* 10(4), 603–616 (2017)
20. Kreisel, G.: Mathematical logic. In: Saaty, T.L. (ed.) *Lectures on modern mathematics*, vol. 3. Wiley (1963)
21. Kunčar, O., Popescu, A.: A Consistent Foundation for Isabelle/HOL. In: *ITP*. pp. 234–252 (2015)
22. Kunčar, O., Popescu, A.: Comprehending Isabelle/HOL’s consistency. In: *ESOP*. pp. 724–749 (2017)
23. Löb, M.: Solution of a Problem of Leon Henkin. *The Journal of Symbolic Logic* 20(2), 115–118 (1955)

24. Nipkow, T., Paulson, L., Wenzel, M.: Isabelle/HOL — A Proof Assistant for Higher-Order Logic, LNCS, vol. 2283. Springer (2002)
25. O'Connor, R.: Essential incompleteness of arithmetic verified by Coq. In: TPHOLs. pp. 245–260 (2005)
26. Paulson, L.C.: A machine-assisted proof of Gödel's incompleteness theorems for the theory of hereditarily finite sets. *Rew. Symb. Logic* 7(3), 484–498 (2014)
27. Paulson, L.C.: A mechanised proof of Gödel's incompleteness theorems using Nominal Isabelle. *J. Autom. Reasoning* 55(1), 1–37 (2015)
28. Paulson, L.C., Blanchette, J.C.: Three years of experience with Sledgehammer, a practical link between automatic and interactive theorem provers. In: The 8th International Workshop on the Implementation of Logics, IWIL 2010, Yogyakarta, Indonesia, October 9, 2011. pp. 1–11 (2010)
29. Popescu, A., Roşu, G.: Term-generic logic. *Theor. Comput. Sci.* 577, 1–24 (2015)
30. Popescu, A., Traytel, D.: A Formally Verified Abstract Account of Gödel's Incompleteness Theorems (Extended Report) (2019), https://bitbucket.org/traytel/abstract_incompleteness/downloads/report.pdf
31. Popescu, A., Traytel, D.: Formalization associated with this paper. https://bitbucket.org/traytel/abstract_incompleteness/ (2019)
32. Quäife, A.: Automated proofs of Löb's theorem and Gödel's two incompleteness theorems. *J. Autom. Reasoning* 4(2), 219–231 (1988)
33. Raatikainen, P.: Gödel's incompleteness theorems. In: The Stanford Encyclopedia of Philosophy. Metaphysics Research Lab, Stanford University (2018)
34. Schlichtkrull, A., Blanchette, J.C., Traytel, D., Waldmann, U.: Formalizing Bachmair and Ganzinger's ordered resolution prover. In: Galmiche, D., Schulz, S., Sebastiani, R. (eds.) *IJCAR 2018*. LNCS, vol. 10900, pp. 89–107. Springer (2018)
35. Shankar, N.: *Metamathematics, Machines, and Gödel's Proof*. Cambridge University Press (1994)
36. Sieg, W.: *Elementary proof theory*. Tech. rep., Institute for Mathematical Studies in the Social Sciences, Stanford (1978)
37. Sieg, W., Field, C.: Automated search for Gödel's proofs. *Ann. Pure Appl. Logic* 133(1-3), 319–338 (2005)
38. Smith, P.: *An introduction to Gödel's incompleteness theorems*. Cambridge University Press (2007)
39. Smoryński, C.: The incompleteness theorems. In: Barwise, J. (ed.) *Handbook of Mathematical Logic*, pp. 821–865. North-Holland (1977)
40. Świerczkowski, S.: Finite sets and Gödel's incompleteness theorems. *Dissertationes Mathematicae* 422, 1–58 (2003)
41. Tarski, A., Mostowski, A., Robinson, R.: *Undecidable Theories*. Studies in Logic and the Foundations of Mathematics. North-Holland (1953), 3rd edition, 1971